# Local linear estimation of a smooth distribution based on censored data

**Liang Peng[1] and Shan Sun[2]**

September 2003

## Abstract

We propose a local linear estimator of a smooth distribution function based on censored data. This new estimator applies local linear techniques to observations from a regression model where the value of the product limit estimator equals the value of the true distribution plus an error term. We show that for most commonly used kernel functions, our local linear estimator has a smaller mean squared error than kernel estimator studied by Ghorai and Susarla (1990).

**Keywords.** Censored data, distribution function, local linear estimation, MSE, Product-limit estimator.

---

[1]School of Mathematics, Georgia Institute of Technology, Atlanta GA 30332-0160, USA. Email: peng@math.gatech.edu

[2]Department of Mathematics and Statistics, Texas Tech University, Lubbock, Texas 79409 - 1042, USA. Email: ssun@math.ttu.edu. Research was supported by NSA Grant MDA904-02-1-0071

# 1   Introduction

Let $X_1, \cdots, X_n$ be i.i.d. random variables with a smooth distribution function $F$. In many practical situations, not all $X_i$'s may be observable. This results in incomplete observations. A common cause for incomplete data is the right censoring, in which case we observe $\{Z_i, \delta_i\}, i = 1, ..., n$, where $Z_i = min(X_i, Y_i)$ and $\delta_i = I(X_i \leq Y_i)$. Thus $\delta_i = 1$ indicates the survival time $X_i$ for the $i$th individual is observed while $\delta_i = 0$ indicates $X_i$ is not observed but it is known to be greater than $Y_i$. Assume that $Y_1, ..., Y_n$ are i.i.d. from the distribution function $G$ and that $X_i$'s and $Y_i$'s are independent. Then the c.d.f. of $Z_i, i = 1, ..., n$ is given by $1 - (1 - F(x))(1 - G(x))$.

Censored data is widely seen in medical studies. Also, a fundamental problem of interest in nonparametric statistical estimation and hypothesis testing is making inference about a c.d.f. without specifying a particular parametric form for $F$.

In estimating the distribution function $F$, a popular nonparametric estimator $F_n$ based on the right censored data $\{Z_i, \delta_i\}, i = 1, ..., n$ is the well-known Kaplan-Meier (Kaplan and Meier, 1958) estimator given by

$$F_n(x) = \begin{cases} 1 - \Pi_{j=1}^n \left[ \frac{N(Z_j)}{1+N(Z_j)} \right]^{I(Z_j \leq x, \delta_j = 1)} & \text{if} \quad x < \max(Z_1, \cdots, Z_n) \\ 1 & \text{elsewhere,} \end{cases}$$

where $N(u) = \sum_{j=1}^n I(Z_j > u)$. The large sample properties of the product-limit estimator have drawn much attention in the literature. See Chen and Lo (1997) and references cited therein.

However, in situations where it is known or it is reasonable to assume that the underlying distribution function $F$ is smooth with density $f$, it is more natural to use a smooth random function as an estimator of $F$. In case of complete data, it was pointed out by Read (1972) that the choice of the empirical function $F_n$ (as an estimator for $F$) does not always lead to the best estimator of $F$ due to the fact that $F_n$ is in admissible with respect to the integrated squared loss. An intuitively appealing and easily understood competitor to $F_n$ is the smooth empirical distribution function proposed by Nadaraya (1964) as follows:

$$\bar{F}_n(x) = h^{-1} \int_{-\infty}^{\infty} k(\frac{x-y}{h}) F_n(y)\, dy, \tag{1.1}$$

where $k$ is a kernel function and $h = h(n) > 0$ is a sequence of bandwidth such that $h \to 0$ as $n \to \infty$.

There is an extensive literature on the study of $\bar{F}_n$, see Reiss (1981), Azzalini (1981) and Falk (1983), among others. A notable property, namely relative deficiency of the empirical estimator with respect to the kernel smooth estimator $\bar{F}_n$ was studied by Falk (1984). He proved that asymptotic performance of $\bar{F}_n$ is better than that of the empirical estimator in terms of second order efficiency.

Kernel smooth estimation for $F$ was extended to the case of the censored data by Ghorai and Susarla (1990). The estimator is of the same form as in (1.1) with $F_n$ being the Kaplan-Meier type of estimator, and kernel $k$ and bandwidth $h$ are the same as in (1.1). Ghoral and Susarla (1990) showed that the kernel estimator is better than the product-limit estimator in terms of mean squared errors. They also provided the weak convergency of the normalized processes $\bar{F}_n(x)$ and strong consistency of $\bar{F}_n$. Later, Sun, Sun and Diao (2001) considered the same estimator $\bar{F}_n$ as well as the corresponding quantile processes $\bar{Q}_n = \bar{F}_n^{-1}$ and obtained the weak convergence of the normalized processes $\bar{F}_n(x)$ and $\bar{Q}_n = \bar{F}_n^{-1}$ in a more general set of conditions.

In this paper we are interested in another type of smooth estimation of $F$ based on censored data. We compare the asymptotic performances of our estimator with the well known Nadaraya type (Ghoral and Susarla (1990)) of kernel estimator in terms of high order mean squared error.

First, we propose our smooth estimator of the underlying distribution function $F$ as follows: observe the following relation

$$F_n(Z_j) = F(Z_j) + \text{error},$$

for those $j's$ such that $\delta_j = 1$, i.e., $Z_j = X_j$. We could apply the local smoothing techniques (see Fan and Gijbels (1996)) to estimating function $F$. Here we concentrate on the local linear estimation. Let $(\hat{a}, \hat{b})$ be the value of $(a, b)$ that minimizes the following kernel weighted squared

errors:

$$\sum_{j=1}^{n}\{F_n(Z_j) - a - b(x - Z_j)\}^2 k(\frac{x - Z_j}{h})\delta_j.$$

Then our local linear estimation is defined as $\hat{a}$ and has the following explicit expression

$$\hat{a} = \tilde{F}_n(x) = \frac{\sum_{j=1}^{n} w_j F_n(Z_j)}{\sum_{j=1}^{n} w_j},$$

where $w_j = \delta_j k(\frac{x-Z_j}{h})[s_{n,2} - (x - Z_j)s_{n,1}], j = 1, ...n$ and $s_{n,l} = \sum_{j=1}^{n} \delta_j k(\frac{x-Z_j}{h})(x - Z_j)^l, l = 1, 2.$

In order to calculate the MSE we work with the modified local linear estimator

$$\hat{F}_n(x) = [o \vee \tilde{F}_n(x)] \wedge 1.$$

We organize this paper as follows. In section 2, we state the main result: mean square error of $\hat{F}_n(x)$. In section 3, we present the simulation studies. All proofs are deferred to section 4.

# 2   Main result

Throughout this paper we assume the following regularity conditions:

(A1) $k$ is a symmetric density with support $(-1, 1)$ and $K(x) = \int_{-1}^{x} k(y)\, dy$;

(A2) $h = h(n) > 0$ is a bandwidth satisfying $d_0 n^{-1+\epsilon_0} \le h \le d_1 n^{-\epsilon_1}$ for some positive numbers of $d_0, d_1, \epsilon_0 \in (\frac{1}{2}, \frac{2}{3}]$ and $\epsilon_1 \in (0, \frac{1}{3}]$;

(A3) $\delta_0 > 0$ is small enough so that it satisfies $1 - 2\epsilon_1 - 2\epsilon_1\delta_0 > 0$ and $-1 - 2\delta_0 + 2\epsilon_0 + 2\epsilon_0\delta_0 > 0$, which implies $\sqrt{nh}h^{1/2+\delta_0} \to \infty$.

(A4) $F''(x)$ and $G'(x)$ are continuous.

Denote

$$c_1 = \int_{-1}^{1} x^2 k(x)\, dx, \quad c_2 = \int_{-1}^{1} xk(x)K(x)\, dx, \quad c_3 = \int_{-1}^{1} x^2 k^2(x)\, dx,$$

4

$H_1(u) = P(Z_j \leq u, \delta_j = 1)$, $H_2(u) = P(Z_j \leq u, \delta_j = 0)$, $H(u) = H_1(u) + H_2(u)$,
$\tau_H = \sup\{t : H(t) < 1\}$ and $\gamma(s) = \int_{-\infty}^{s}[1 - G(u)]^{-1}[1 - F(u)]^{-2} \, dF(u)$.

Then our main result is the following.

**Theorem 1.** *Assume* **(A1)** *-* **(A4)** *are true and* $x < \tau_H$. *Then*

$$
\begin{aligned}
&MSE(\hat{F}_n(x)) \\
=\ & \gamma(x)[H_1'(x)]^2[1 - F(x)]^2 n^{-1} \\
& -[4c_2 - c_3][1 - G(x)][F'(x)]^3 n^{-1} h \\
& +\tfrac{1}{4}c_1^2[F''(x)]^2 h^4 + o(h^4 + n^{-1}h)
\end{aligned}
\tag{2.1}
$$

Our goal here is to compare the asymptotic performance of our local linear estimator $\hat{F}_n$ with that of the kernel estimator $\bar{F}_n$ in terms of MSE. Ghorai and Susarla (1990) establish the MSE of $\bar{F}_n$ using estimated bias term. For comparison purpose, we can prove, under the same set of conditions, the exact MSE of $\bar{F}_n$ is given by following:

$$
\begin{aligned}
&MSE(\bar{F}_n(x)) \\
=\ & \gamma(x)[H_1'(x)]^2[1 - F(x)]^2 n^{-1} \\
& -2c_2[1 - G(x)][F'(x)]^3 n^{-1} h \\
& +\tfrac{1}{4}c_1^2[F''(x)]^2 h^4 + o(h^4 + n^{-1}h)
\end{aligned}
\tag{2.2}
$$

The proof of (2.2) is similar to the proof of our main theorem, hence is omitted here.

Observing (2.1) and (2.2), we notice that the main difference between two equations is the coefficients in the first order of the bandwidth $h$. It was pointed out by Cheng and Peng (2002), $2c_2 - c_3$ is positive for most conventional kernel functions. See Table 1 below. Therefore our local linear estimator $\hat{F}_n(x)$ has a smaller mean squared error than kernel estimator $\bar{F}_n(x)$ for most commonly used kernels.

**Table 1**: *Values of $2c_2 - c_3$ for some commonly used kernels.*

| Kernel | $2c_2 - c_3$ |
|---|---|
| Epanechnikov $k(x) = \frac{3}{4}(1 - x^2)I(|x| \leq 1)$ | $\frac{6}{35}$ |
| Biweight $k(x) = \frac{15}{16}(1 - x^2)^2 I(|x| \leq 1)$ | $\frac{5}{33}$ |
| Triangular $k(x) = (1 - |x|)I(|x| \leq 1)$ | $\frac{1}{6}$ |
| Uniform $k(x) = \frac{1}{2}I(|x| \leq 1)$ | $\frac{1}{6}$ |

Since the bandwidth plays a critical role in implementing practical estimation and it determines the trade-off between the amount of smoothness obtained and closeness of the estimation to the true distribution, it is important that we provide a way to select the optimal bandwidth. As in the cases of smooth distribution function estimation (see Altman and Leger (1995), Bowman, Hall and Prvan (1998) or Cheng and Peng (2002)) we choose optimal bandwidth $h$ in the sense of minimizing the second order term in the expansion of $MSE(\hat{F}_n(x))$, i.e., minimizing

$$-[4c_2 - c_3][1 - G(x)][F'(x)]^3 n^{-1}h + \frac{1}{4}c_1^2[F''(x)]^2 h^4 + o(h^4 + n^{-1}h),$$

which gives

$$h_l^* = \{\frac{[4c_2 - c_3][1 - G(x)][F'(x)]^3}{c_1^2[F''(x)]^2}\}^{1/3}n^{-1/3}.$$

Similarly the optimal bandwidth for kernel smooth estimator $\bar{F}_n(x)$ can be obtained as

$$h_k^* = \{\frac{2c_2[1 - G(x)][F'(x)]^3}{c_1^2[F''(x)]^2}\}^{1/3}n^{-1/3}.$$

## 3   Simulation study

A Monte Carlo study was conducted to compare the performance between the local linear estimator $\hat{F}_n$ and kernel smooth distribution estimator $\bar{F}_n$ in terms of mean squared error. The Epanechnikov kernel defined by $k(x) = \frac{3}{4}(1 - x^2)I(|x| \leq 1)$ was used to construct the two estimators.

Let $F(x) = 1 - e^{-x}$ and $G(x) = 1 - e^{-\alpha x}$, where $\alpha$ is the censoring parameter. The relative efficiencies of $\hat{F}_n(x)$ and $\bar{F}_n(x)$ with respect to $F_n(x)$ ( ratios of the mean squared errors of $\hat{F}_n(x)$ and $\bar{F}_n(x)$ to that of $F_n(x)$, respectively) were computed and plotted against different $h$ (for $h = 0.01$ with the increment of 0.01) for $x = -\log(0.8), -\log(0.6), -\log(0.4)$ and $-\log(0.2)$. In Figures 1-4, $\alpha$ is chosen to have 30% censoring ($\alpha = 3/7$) while in figures 5-8 $\alpha$ is chosen to have 5% censoring ($\alpha = 1/19$). The values of the MSE were the averages based on 1000 repetitions of sample of size 200. From the figures 1-8, we observe the following facts:

1. Both kernel and local linear estimators have smaller MSE than the product-limit estimator. In general, the plots show that the optimal mean squared error of our local linear estimator is smaller than that of kernel smooth estimator. Also, the local linear estimator is better than kernel estimator when both estimators are based on their optimal bandwidths (that is, the bandwidth which minimize the MSE).

2. When $x = -\log(0.6)$ and $x = -\log(0.4)$, that is when $F(x)$ is in the neighborhood of 0.5, our local linear estimator performs much better than kernel estimator (see Figures 2, 3, 6 and 7).

3. For small $h$, there are some problems with using local linear estimator when $F(x)$ is near zero or one, this is due to the fact that the local linear estimator is not an increasing function. On the other hand, when $F(x)$ is close to zero and one, our local linear estimator is quite robust against the bandwidth $h$. computed for x.$h_l^* = 0.258, h_k^* = 0.217$ $h_l^* = 0.162, h_k^* = 0.137$
$h_l^* = 0.307, h_k^* = 0.259$ $h_l^* = 0.221, h_k^* = 0.186$
$h_l^* = 0.225, h_k^* = 0.190$ $h_l^* = 0.142, h_k^* = 0.120$
$h_l^* = 0.268, h_k^* = 0.226$

# 4 Proofs of Main Theorems

Denote

$$H_{n,1}(u) = \frac{1}{n} \sum_{j=1}^{n} I(Z_j \leq u, \delta_j = 1),$$

$$H_{n,2}(u) = \frac{1}{n} \sum_{j=1}^{n} I(Z_j \le u, \delta_j = 0),$$

$$H_n(u) = H_{n,1}(u) + H_{n,2}(u),$$

and also define the following sets:

$$A_1 = \{|n^{-2}h^{-4} \sum_{j=1}^{n} w_j - (H_1'(x))^2 c_1| > h^{1/2+\delta_0}\}$$

$$A_2 = \{|n^{-1}h^{-3}s_{n,1} + H_1''(x)c_1| > h^{1/2+\delta_0}\}$$

$$A_3 = \{|n^{-1}h^{-3}s_{n,2} - H_1'(x)c_1| > h^{1/2+\delta_0}\}.$$

By applying Edgeworth expansions, we obtain following results:

$$\begin{cases} P(A_1) = O(n^{-3}h^{-3}) \\ P(A_2) = O(n^{-3}h^{-3}) \\ P(A_3) = O(n^{-3}h^{-3}) \end{cases} \tag{4.1}$$

For proofs of (4.1), see Hall (1992).

Let $T$ be such that $1 - H(T) > d$ with some $d > 0$ and $M, \lambda$ denote generic positive constants. Then it follows from Major and Rejto (1988) that the process $\{F_n(u) - F(u), -\infty < u < \infty, 1 - H(u) > 0\}$ can be represented as

$$F_n(u) - F(u) = (1 - F(u))[B_1(n,u) + B_2(n,u)] + R(n,u),$$

where

$$B_1(n,u) = \frac{H_{n,1}(u) - H_1(u)}{1 - H(u)} - \int_{-\infty}^{u} \frac{H_{n,1}(y) - H_1(y)}{[1 - H(y)]^2} dH(y),$$

$$B_2(n,u) = \int_{-\infty}^{u} \frac{H_n(y) - H(y)}{[1 - H(y)]^2} dH_2(y),$$

and

$$P(A_4) \le M e^{-\lambda h^{-\delta_0}},$$

where

$$A_4 = \{\sup_{u \le T} n|R(n,u)| > h^{-\delta_0}\}.$$

Moreover there exists a Gaussian Process $W(u), -\infty < u < \infty$, with $E(W(u)) = 0$ and covariance

$$E(W(s)W(t)) = \gamma(s) \qquad (4.2)$$

for $-\infty < s \le t < \infty$, where $\gamma(s)$ is defined right before Theorem 1, which satisfies

$$\begin{cases} P(A_5) \le Me^{-\lambda h^{-\delta_0}} \\ P(A_6) \le Me^{-\lambda h^{-\delta_0}}, \end{cases} \qquad (4.3)$$

where

$$A_5 = \{ \sup_{-\infty < u \le T} \sqrt{n} |\sqrt{n}[B_1(n, u) + B_2(n, u)] - W(u)| > h^{-\delta_0} \}$$

and

$$A_6 = \{ \sup_{-\infty < u < \infty} \sqrt{n} |H_n(u) - H(u)| > h^{-\delta_0} \}.$$

Using (4.1) - (4.3) and note that $H_1'(x) = [1 - G(x)]F'(x)$ and $\gamma'(x) = [1 - G(x)]^{-1}[1 - F(x)]^{-2}F'(x)$, it is clear that to prove Theorem 1, we need to prove the following:

$$E\{[\hat{F}_n(x) - F(x)]^2 I(\cap_{j=1}^6 A_j^c)\}$$

$$= c_1^2 \gamma(x)[H_1'(x)]^4[1 - F(x)]^2 n^{-1} - 2c_1^2 c_2 \gamma'(x)[H_1'(x)]^4[1 - F(x)]^2 n^{-1}h$$

$$-c_1^2[2c_2 - c_3]\gamma'(x)[H_1'(x)]^2[1 - H(x)]^2[F'(x)]^2 n^{-1}h + \tfrac{1}{4}c_1^4[H_1'(x)]^4[F''(x)]^2 h^4 + o(n^{-1}h + h^4), \qquad (4.4)$$

where $A^c$ denotes the complementary set of $A$ and $I(A)$ denotes the indicator function of set $A$.

We first show that (4.4) is true when $\hat{F}_n(x)$ is replaced by $\tilde{F}_n(x)$. Therefore similarly we can prove that (4.4) holds. To this end, we express the term $n^{-2}h^{-4}\sum_{j=1}^n w_j[F_n(Z_j) - F(x)]$ as follows:

9

$$n^{-2}h^{-4}\sum_{j=1}^{n} w_j[F_n(Z_j) - F(x)]$$

$$= [n^{-1}h^{-3}s_{n,2} - H_1'(x)c_1]\{n^{-1}h^{-1}\sum_{j=1}^{n}\delta_j k(\tfrac{x-Z_j}{h})[F_n(Z_j) - F(x)]\}$$

$$-[n^{-1}h^{-3}s_{n,1} + H_1''(x)c_1]\{n^{-1}h^{-1}\sum_{j=1}^{n}(x - Z_j)\delta_j k(racx - Z_j h)[F_n(Z_j) - F(x)]\}$$

$$+c_1 n^{-1}h^{-1}\sum_{j=1}^{n}\delta_j k(\tfrac{x-Z_j}{h})[H_1'(x) + (x - Z_j)H_1''(x)][F(Z_j) - F(x)]$$

$$+c_1 n^{-1}h^{-1}\sum_{j=1}^{n}\delta_j k(\tfrac{x-Z_j}{h})[H_1'(x) + (x - Z_j)H_1''(x)][F_n(Z_j) - F(Z_j)],$$

$$= I_1 + I_2 + I_3 + I_4$$

(4.5)

Note that

$$dH_{n,1}(u) = [1 - H(u)]\,d[B_1(n, u) + B_2(n, u)] + H_1'(u)\,du - \frac{H_n(u) - H(u)}{1 - H(u)}H_2'(u)\,du,$$

therefore, $I_3$ and $I_4$ can be further decomposed as follows:

$$
\begin{aligned}
I_3 &= c_1 h^{-1} \int_{-\infty}^{\infty} k(\tfrac{x-s}{h})[H_1'(x) + (x-s)H_1''(x)][F(s) - F(x)]\, dH_{n,1}(s) \\
&= c_1 h^{-1} \int_{-\infty}^{\infty} k(\tfrac{x-s}{h})[H_1'(x) + (x-s)H_1''(x)][F(s) - F(x)][1 - H(s)] \\
&\quad \times d[B_1(n,s) + B_2(n,s)] \\
&\quad + c_1 h^{-1} \int_{-\infty}^{\infty} k(\tfrac{x-s}{h})[H_1'(x) + (x-s)H_1''(x)][F(s) - F(x)]H_1'(s)\, ds \\
&\quad - c_1 h^{-1} \int_{-\infty}^{\infty} k(\tfrac{x-s}{h})[H_1'(x) + (x-s)H_1''(x)][F(s) - F(x)]\tfrac{H_n(s) - H(s)}{1 - H(s)}H_2'(s)\, ds \\
&= c_1 h^{-1} n^{-1/2} \int_{-1}^{1} \{\sqrt{n}[B_1(n, x-sh) + B_2(n, x-sh)] - W(x-sh)\} \\
&\quad \times d\{k(s)[H_1'(x) + shH_1''(x)][F(x-sh) - F(x)][1 - H(x-sh)]\} \\
&\quad + c_1 h^{-1} n^{-1/2} \int_{-1}^{1} W(x-sh)\, d\{k(s)[H_1'(x) + shH_1''(x)] \\
&\quad \times [F(x-sh) - F(x)][1 - H(x-sh)]\} \\
&\quad + c_1 \int_{-1}^{1} k(s)[H_1'(x) + shH_1''(x)][F(x-sh) - F(x)]H_1'(x-sh)\, ds \\
&\quad - c_1 \int_{-1}^{1} k(s)[H_1'(x) + shH_1''(x)][F(x-sh) - F(x)]\tfrac{H_n(x-sh) - H(x-sh)}{1 - H(x-sh)}H_2'(x-sh)\, ds \\
&= II_1 + II_2 + II_3 + II_4
\end{aligned}
$$

$$(4.6)$$

and

$$
\begin{aligned}
I_4 \;=\;& c_1 h^{-1} \int_{-\infty}^{\infty} k\left(\tfrac{x-s}{h}\right)[H_1'(x) + (x-s)H_1'(x)][F_n(s) - F(s)]\, dH_{n,1}(s) \\[2mm]
=\;& c_1 n^{-1} h^{-1} \sum_{j=1}^{n} \delta_j k\left(\tfrac{x-Z_j}{h}\right)[H_1'(x) + (x-Z_j)H_1''(x)]R(n, Z_j) \\[2mm]
& + c_1 h^{-1} \int_{-\infty}^{\infty} k\left(\tfrac{x-s}{h}\right)[H_1'(x) + (x-s)H_1''(x)][1 - F(s)][B_1(n,s) + B_2(n,s)]\, dH_{n,1}(s) \\[2mm]
=\;& c_1 n^{-1} h^{-1} \sum_{j=1}^{n} \delta_j k\left(\tfrac{x-Z_j}{h}\right)[H_1'(x) + (x-Z_j)H_1''(x)]R(n, Z_j) \\[2mm]
& + \tfrac{1}{2} c_1 h^{-1} \int_{-\infty}^{\infty} k\left(\tfrac{x-s}{h}\right)[H_1'(x) + (x-s)H_1''(x)][1 - F(s)][1 - H(s)] \times \, d[B_1(n,s) + B_2(n,s)]^2 \\[2mm]
& + c_1 h^{-1} \int_{-\infty}^{\infty} k\left(\tfrac{x-s}{h}\right)[H_1'(x) + (x-s)H_1''(x)][1 - F(s)][B_1(n,s) + B_2(n,s)]H_1'(s)\, ds \\[2mm]
& - c_1 h^{-1} \int_{-\infty}^{\infty} k\left(\tfrac{x-s}{h}\right)[H_1'(x) + (x-s)H_1''(x)][1 - F(s)][B_1(n,s) + B_2(n,s)] \\[2mm]
& \times \frac{H_n(s) - H(s)}{1 - H(s)} H_2'(s)\, ds \\[2mm]
=\;& c_1 n^{-1} h^{-1} \sum_{j=1}^{n} \delta_j k\left(\tfrac{x-Z_j}{h}\right)[H_1'(x) + (x-Z_j)H_1''(x)]R(n, Z_j) \\[2mm]
& + \tfrac{1}{2} c_1 h^{-1} \int_{-1}^{1} [B_1(n, x - sh) + B_2(n, x - sh)]^2 \\[2mm]
& \times \, d\{k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)][1 - H(x - sh)]\} \\[2mm]
& + c_1 \int_{-1}^{1} k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)][B_1(n, x - sh) + B_2(n, x - sh)] \\[2mm]
& \times H_1'(x - sh)\, ds \\[2mm]
& - c_1 \int_{-1}^{1} k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)][B_1(n, x - sh) + B_2(n, x - sh)] \\[2mm]
& \times \frac{H_n(x - sh) - H(x - sh)}{1 - H(x - sh)} H_2'(x - sh)\, ds
\end{aligned}
$$

$$= c_1 n^{-1} h^{-1} \sum_{j=1}^{n} \delta_j k(\tfrac{x-Z_j}{h})[H_1'(x) + (x - Z_j)H_1''(x)]R(n, Z_j)$$

$$+ \tfrac{1}{2} c_1 n^{-1} h^{-1} \int_{-1}^{1} \{\sqrt{n}[B_1(n, x - sh) + B_2(n, x - sh)] - W(x - sh)\}^2$$

$$\times d\{k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)][1 - H(x - sh)]\}$$

$$+ c_1 n^{-1} h^{-1} \int_{-1}^{1} \{\sqrt{n}[B_1(n, x - sh) + B_2(n, x - sh)] - W(x - sh)\} \times W(x - sh)$$

$$\times d\{k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)][1 - H(x - sh)]\}$$

$$+ \tfrac{1}{2} c_1 n^{-1} h^{-1} \int_{-1}^{1} W^2(x - sh)$$

$$\times d\{k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)][1 - H(x - sh)]\}$$

$$+ c_1 n^{-1/2} \int_{-1}^{1} k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)]$$

$$\times \{\sqrt{n}[B_1(n, x - sh) + B_2(n, x - sh)] - W(x - sh)\}H_1'(x - sh)\, ds$$

$$+ c_1 n^{-1/2} \int_{-1}^{1} k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)]W(x - sh)H_1'(x - sh)\, ds$$

$$- c_1 n^{-1/2} \int_{-1}^{1} k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)]$$

$$\times \{\sqrt{n}[B_1(n, x - sh) + B_2(n, x - sh)] - W(x - sh)\}$$

$$\times \tfrac{H_n(x-sh) - H(x-sh)}{1 - H(x-sh)} H_2'(x - sh)\, ds$$

$$- c_1 n^{-1/2} \int_{-1}^{1} k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)]W(x - sh)$$

$$\times \tfrac{H_n(x-sh) - H(x-sh)}{1 - H(x-sh)} H_2'(x - sh)\, ds \tag{4.7}$$

$$= III_1 + \cdots + III_8.$$

The terms $II_1, II_3$ and $II_4$ in (4.6) can be estimated as follows:

$$\begin{cases} |II_1 I(\cap_{j=1}^{6} A_j^c)| \le M n^{-1} h^{-\delta_0} \\[2mm] |II_4 I(\cap_{j=1}^{6} A_j^c)| \le M n^{-1/2} h^{1-\delta_0} \end{cases} \tag{4.8}$$

and

$$II_3 = c_1 \int_{-1}^{1} k(s)[H_1'(x) + shH_1''(x)][-shF'(x) + \tfrac{1}{2}s^2h^2F''(x) + o(h^2)]$$

$$\times [H_1'(x) - shH_1''(x) + o(h)]\,ds \qquad (4.9)$$

$$= h^2\tfrac{1}{2}c_1^2[H_1'(x)]^2F''(x) + o(h^2).$$

The terms $III_1 - -III_5$, $III_7$ and $III_8$ in (4.7) can be estimated as follows:

$$E\{III_1^2 I(\cap_{j=1}^6 A_j^c)\}$$

$$\leq c_1^2 n^{-2}h^{-2}n^{-2}h^{-2\delta_0} E\{\sum_{j=1}^{n} \delta_j k(\tfrac{x-Z_j}{h})[H_1'(x) + (x - Z_j)H_1''(x)]\}^2 \qquad (4.10)$$

$$\leq Mn^{-2}h^{-2\delta_0},$$

and

$$\begin{cases}
|III_2 I(\cap_{j=1}^6 A_j^c)| \leq Mn^{-2}h^{-1-2\delta_0} \\[2mm]
|III_3 I(\cap_{j=1}^6 A_j^c)| \leq Mn^{-3/2}h^{-1-\delta_0} \\[2mm]
|E(III_4^2 I(\cap_{j=1}^6 A_j^c))| \leq E(III_4^2) \leq Mn^{-2}h^{-1} \\[2mm]
|III_5 I(\cap_{j=1}^6 A_j^c)| \leq Mn^{-1}h^{-\delta_0} \\[2mm]
|III_7 I(\cap_{j=1}^6 A_j^c)| \leq Mn^{-1}h^{-2\delta_0} \\[2mm]
|III_8 I(\cap_{j=1}^6 A_j^c)| \leq Mn^{-1}h^{-\delta_0}.
\end{cases} \qquad (4.11)$$

Now the only terms remained to be analyzed are $II_2$ in (4.6) and $III_6$ in (4.7). We proceed in the following. Put

$$Q_1(s) = k(s)[H_1'(x) + shH_1''(x)][F(x - sh) - F(x)][1 - H(x - sh)],$$

$$Q_2(s) = k(s)[H_1'(x) + shH_1''(x)][1 - F(x - sh)]H_1'(x - sh),$$

$$Q_3(s) = hk(s)\{-sH_1'(x)[1 - H(x)]F'(x) - s^2hH_1''(x)[1 - H(x)]F'(x) -$$

$$-s^2hH_1'(x)H'(x)F'(x) + \tfrac{1}{2}s^2hH_1'(x)[1 - H(x)]F''(x)\},$$

and

$$Q_4(s) = k(s)\{[H_1'(x)]^2[1 - F(x)] + sh[H_1'(x)]^2F'(x)\}.$$

14

Then

$$Q_1(s) = Q_3(s) + o(h^2), \quad Q_1'(s) = Q_3'(s) + o(h^2), \quad Q_2(s) = Q_4(s) + o(h)$$

Observe that the following equations hold:

$$E\{\int_{-1}^1 W(x - sh)Q_3'(s)\,ds\}^2$$

$$= 2\int_{-1}^1\{\int_{-1}^s \gamma(x - sh)Q_3'(s)Q_3'(t)\,dt\}\,ds$$

$$= 2\int_{-1}^1 \gamma(x - sh)Q_3'(s)Q_3(s)\,ds$$

$$= 2\int_{-1}^1\{\gamma(x) - sh\gamma'(x) + o(h)\}Q_3'(s)Q_3(s)\,ds \tag{4.12}$$

$$= 2\int_{-1}^1\{-sh\gamma'(x) + o(h)\}Q_3'(s)Q_3(s)\,ds$$

$$= h^3 c_3 \gamma'(x)[H_1'(x)]^2[1 - H(x)]^2[F'(x)]^2 + o(h^3),$$

$$E\{\int_{-1}^1 W(x - sh)Q_4(s)\,ds\}^2$$

$$= 2\int_{-1}^1\{\int_{-1}^s \gamma(x - sh)Q_4(s)Q_4(t)\,dt\}\,ds$$

$$= 2\int_{-1}^1 \gamma(x - sh)Q_4(s)\{K(s)[H_1'(x)]^2[1 - F(x)] + K(s)sh[H_1'(x)]^2 F'(x) -$$

$$- \int_{-1}^s K(t)h[H_1'(x)]^2 F'(x)\,dt\}\,ds$$

$$= 2\int_{-1}^1 \gamma(x - sh)Q_4(s)K(s)\{[H_1'(x)]^2[1 - F(x)] + sh[H_1'(x)]^2 F'(x)\}\,ds \tag{4.13}$$

$$-2\int_{-1}^1 K(t)h[H_1'(x)]^2 F'(x)\int_t^1 \gamma(x - sh)Q_4(s)\,ds\}\,dt$$

$$= 2\int_{-1}^1 k(s)K(s)\{\gamma(x)[H_1'(x)]^4[1 - F(x)]^2 + 2sh\gamma(x)[H_1'(x)]^4 F'(x)[1 - F(x)] -$$

$$- sh\gamma'(x)[H_1'(x)]^4[1 - F(x)]^2\}\,ds$$

$$-2\int_{-1}^1 K(t)[1 - K(t)]h\gamma(x)[H_1'(x)]^4 F'(x)[1 - F(x)]\,dt + o(h)$$

$$= \gamma(x)[H_1'(x)]^4[1 - F(x)]^2 - h2c_2\gamma'(x)[H_1'(x)]^4[1 - F(x)]^2 + o(h)$$

and

15

$$2E\{\int_{-1}^{1}\int_{-1}^{1}W(x-sh)Q_3'(s)W(x-th)Q_4(t)\,dsdt\}$$

$$= \ 2\int_{-1}^{1}\{\int_{-1}^{s}\gamma(x-sh)Q_3'(s)Q_4(t)\,dt\}\,ds$$

$$+2\int_{-1}^{1}\{\int_{s}^{1}\gamma(x-th)Q_3'(s)Q_4(t)\,dt\}\,ds$$

$$= \ 2\int_{-1}^{1}\{\int_{-1}^{s}\gamma(x-sh)Q_3'(s)Q_4(t)\,dt\}\,ds$$

$$+2\int_{-1}^{1}\{\int_{-1}^{t}\gamma(x-th)Q_3'(s)Q_4(t)\,ds\}\,dt$$

$$= \ 2\int_{-1}^{1}\gamma(x-sh)Q_3'(s)K(s)[H_1'(x)]^2[1-F(x)]\,ds$$

$$+2\int_{-1}^{1}\gamma(x-sh)Q_3'(s)K(s)sh[H_1'(x)]^2F'(x)\,ds$$

$$-2\int_{-1}^{1}\gamma(x-sh)Q_3'(s)\{\int_{-1}^{s}K(t)h[H_1'(x)]^2F'(x)\,dt\}\,ds$$

$$+2\int_{-1}^{1}\gamma(x-th)Q_3(t)Q_4(t)\,dt$$

$$= \ 2\int_{-1}^{1}\gamma(x)Q_3'(s)K(s)[H_1'(x)]^2[1-F(x)]\,ds$$

$$-2\int_{-1}^{1}\gamma'(x)shQ_3'(s)K(s)[H_1'(x)]^2[1-F(x)]\,ds$$

$$+2\int_{-1}^{1}\gamma(x)Q_3'(s)K(s)sh[H_1'(x)]^2F'(x)\,ds$$

$$-2\int_{-1}^{1}\{\int_{t}^{1}\gamma(x-sh)Q_3'(s)K(t)h[H_1'(x)]^2F'(x)\,ds\}\,dt$$

$$+2\int_{-1}^{1}\gamma(x)Q_3(t)Q_4(t)\,dt$$

$$-2\int_{-1}^{1}\gamma'(x)thQ_4(t)Q_3(t)\,dt+o(h^2)$$

$$= \ -2\int_{-1}^{1}Q_3(s)k(s)\gamma(x)[H_1'(x)]^2[1-F(x)]\,ds$$

$$+2\int_{-1}^{1}Q_(s)[K(s)+sk(s)]h\gamma'(x)[H_1'(x)]^2[1-F(x)]\,ds$$

$$-2\int_{-1}^{1}Q_3(s)[K(s)+sk(s)]h\gamma(x)[H_1'(x)]^2F'(x)\,ds$$

$$+2\int_{-1}^{1}\gamma(x)Q_3(t)K(t)h[H_1'(x)]^2F'(x)\,dt$$

$$+2\int_{-1}^{1}\gamma(x)Q_3(t)Q_4(t)\,dt$$

$$+2c_3h^2\gamma'(x)[H_1'(x)]^3[1-H(x)]F'(x)[1-F(x)]+o(h^2)$$

16

$$
\begin{aligned}
= \ & -2 \int_{-1}^{1} Q_3(s) Q_4(s) \gamma(x) \, ds \\[4pt]
& +2 \int_{-1}^{1} Q_3(s) k(s) s h [H_1'(x)]^2 F'(x) \gamma(x) \, ds \\[4pt]
& -2 \int_{-1}^{1} h^2 k(s) s [K(s) + s k(s)] [H_1'(x)]^3 [1 - H(x)] F'(x) [1 - F(x)] \gamma'(x) \, ds \\[4pt]
& +2 \int_{-1}^{1} h^2 k(s) s [K(s) + s k(s)] [H_1'(x)]^3 [1 - H(x)] [F'(x)]^2 \gamma(x) \, ds \\[4pt]
& -2 \int_{-1}^{1} h^2 \gamma(x) t k(t) K(t) [H_1'(x)]^3 [1 - H(x)] [F'(x)]^2 \, ds \\[4pt]
& +2 \int_{-1}^{1} Q_3(t) Q_4(t) \gamma(x) \, dt \\[4pt]
& +2 c_3 h^2 \gamma'(x) [H_1'(x)]^3 [1 - H(x)] F'(x) [1 - F(x)] + o(h^2) \\[6pt]
= \ & -h^2 2 c_3 [H_1'(x)]^3 [1 - H(x)] [F'(x)]^2 \gamma(x) \\[4pt]
& -h^2 2 [c_2 + c_3] [H_1'(x)]^3 [1 - H(x)] F'(x) [1 - F(x)] \gamma'(x) \\[4pt]
& +h^2 2 [c_2 + c_3] [H_1'(x)]^3 [1 - H(x)] [F'(x)]^2 \gamma(x) \\[4pt]
& -h^2 2 c_2 [H_1'(x)]^3 [1 - H(x)] [F'(x)]^2 \gamma(x) \\[4pt]
& +h^2 2 c_3 [H_1'(x)]^3 [1 - H(x)] F'(x) [1 - F(x)] \gamma'(x) + o(h^2) \\[6pt]
= \ & -h^2 2 c_2 [H_1'(x)]^3 [1 - H(x)] F'(x) [1 - F(x)] \gamma'(x) + o(h^2).
\end{aligned}
\tag{4.14}
$$

Therefore by (4.12) - (4.14) we have

$$E\{II_2 + III_6\}^2$$

$$= n^{-1}c_1^2\gamma(x)[H_1'(x)]^4[1 - F(x)]^2$$

$$-n^{-1}h2c_1^2c_2\gamma'(x)[H_1'(x)]^4[1 - F(x)]^2$$

$$+n^{-1}hc_1^2c_3\gamma'(x)[H_1'(x)]^2[1 - H(x)]^2[F'(x)]^2$$

$$-n^{-1}h2c_1^2c_2[H_1'(x)]^3[1 - H(x)]F'(x)[1 - F(x)]\gamma'(x) + o(n^{-1}h)$$

$$= n^{-1}c_1^2\gamma(x)[H_1'(x)]^4[1 - F(x)]^2$$

$$-n^{-1}h2c_1^2c_2\gamma'(x)[H_1'(x)]^4[1 - F(x)]^2$$

$$-n^{-1}hc_1^2[2c_2 - c_3]\gamma'(x)[H_1'(x)]^2[1 - H(x)]^2[F'(x)]^2 + o(n^{-1}h).$$

Further

$$E\{(II_2 + III_6)^2 I(\cap_{j=1}^6 A_j^c)\}$$

$$= E\{II_2 + III_6\}^2 - E\{(II_2 + III_6)^2 I(\cup_{j=1}^6 A_j)\} \qquad (4.15)$$

$$= E\{II_2 + III_6\}^2 + o(n^{-1}h).$$

Finally, the terms $I_1, ..., I_4$ in (4.5) can be estimated by using (4.8)–(4.11) and (4.15), that is

$$|E\{I_i I_j I(\cap_{j=1}^6 A_j^c)\}| = o(n^{-1}h + h^4) \qquad (4.16)$$

for $i \neq j$. Thus it follows from (4.8)–(4.11) and (4.16) that (4.4) is true by replacing $\hat{F}_n(x)$ by $\tilde{F}_n(x)$. Similarly we can show that for any $q > 1$

$$E\{(\tilde{F}_n(x) - F(x))^2 I(\cap_{j=1}^6 A_j^c)\}^q = O(n^{-q}). \qquad (4.17)$$

Hence, for any $q_1^{-1} + q_2^{-1} = 1$ with $q_1 > 1$

$$E\{(\tilde{F}_n(x) - F(x))^2 I(\tilde{F}_n(x) < 0) I(\cap_{j=1}^6 A_j^c)\}$$

$$\leq \{E((\tilde{F}_n(x) - F(x))^2 I(\cap_{j=1}^6 A_j^c))^{q_1}\}^{1/q_1} \{EI(\tilde{F}_n(x) < 0) I(\cap_{j=1}^6 A_j^c)\}^{1/q_2}$$

$$\leq \{E((\tilde{F}_n(x) - F(x))^2 I(\cap_{j=1}^6 A_j^c))^{q_1}\}^{1/q_1} \{EI(|\tilde{F}_n(x) - F(x)| > F(x)) I(\cap_{j=1}^6 A_j^c)\}^{1/q_2}$$

$$= \{E((\tilde{F}_n(x) - F(x))^2 I(\cap_{j=1}^6 A_j^c))^{q_1}\}^{1/q_1} \{EI(|\tilde{F}_n(x) - F(x)| I(\cap_{j=1}^6 A_j^c) > F(x))\}^{1/q_2}$$

$$\leq \{E((\tilde{F}_n(x) - F(x))^2 I(\cap_{j=1}^6 A_j^c))^{q_1}\}^{1/q_1} \{\frac{E((\tilde{F}_n(x) - F(x))^2 I(\cap_{j=1}^6 A_j^c))}{F^2(x)}\}^{1/q_2}$$

$$= O(n^{-1-1/q_2})$$

$$= o(n^{-1}h + h^4) \quad (\text{letting } q_2 \to 1).$$

Similarly

$$E\{(\tilde{F}_n(x) - F(x))^2 I(\tilde{F}_n(x) > 1) I(\cap_{j=1}^6 A_j^c)\} = o(n^{-1}h + h^4).$$

Thus, (4.4) holds by noting that

$$E\{(\hat{F}_n(x) - F(x))^2 I(\cap_{j=1}^6 A_j^c)\}$$

$$= E\{(\tilde{F}_n(x) - F(x))^2 I(\cap_{j=1}^6 A_j^c)\}$$

$$-E\{(\tilde{F}_n(x) - F(x))^2 I(\tilde{F}_n(x) < 0) I(\cap_{j=1}^6 A_j^c)\}$$

$$-E\{(\tilde{F}_n(x) - F(x))^2 I(\tilde{F}_n(x) > 1) I(\cap_{j=1}^6 A_j^c)\}.$$

Hence Theorem 1.

# References

[1] N. Altman and C. Leger (1995). Bandwidth selection for kernel distribution function estimation. *J. Statist. Plan. Inf., 46, 195 - 214.*

[2] A. Azzalini (1981). A note on estimation of a distribution function and quantiles by a kernel method. *Biometrika, 68, 326 - 328.*

[3] N. Bowman, P. Hall and T. Prvan (1998). Bandwidth selection for the smoothing of distribution functions. *Biometrika, 85(4), 799 - 808.*

[4] K. Chen and S. Lo (1997). On the rate of uniform convergence of the product-limit estimator: strong and weak laws. *Ann. Statist., 25(3), 1050 1087.*

[5] M. Cheng and L. Peng (2002). Regression modeling for nonparametric estimation of distribution and quantile functions. *Statistica Sinica, 12(4).*

[6] M. Falk (1983). Relative efficiency and deficiency of kernel type estimators of smooth distribution

[7] M. Falk (1984). Relative efficiency of kernel type estimators of quatiles. *Ann. Statist. 12, 261-268*

[8] J. Fan and I. Gijbels (1996). Local Polynomial Modelling and its Applications.*Chapman and Hall, London.*

[9] J.K. Ghorai and V. Susarla (1990). Kernel estimation of a smooth distributon function based on censored data. *Metrika, 37, 71 - 86.*

[10] P. Hall (1992). The Bootstrap and Edgeworth Expansion. *Springer.*

[11] E.L. Kaplan and P. Meier (1958). Nonparametric estimation from incomplete observations. *J. Amer. Statist. Assoc., 53, 457 - 481.*

[12] P. Major and L. Rejto (1988). Strong embedding of the estimator of the disribution function under random censorship. *Ann. Statist., 16, 1113 - 1132.*

[13] E.A. Nadaraya (1964) Some new estimates for distribution functions. *Theory Probab. Appl. 9, 497-500.*

[14] R.R. Read (1972). The asymptotic inadmissibility of the sample distribution function. *Ann. Math. Statist., 43, 89-95.*

[15] R.D. Reiss (1981). Nonparametric estimation of smooth distribution functions. *Scand. J. Statist., 8, 116 - 119.*

kernel quantile estimators

[16] Y. Sun, S. Sun and Y. Diao (2001). Smooth quantile processes from right censored data and construction of simultaneous conference bands. *Communication in Statistics, Theory and Methods. 30 no.4 707-727*
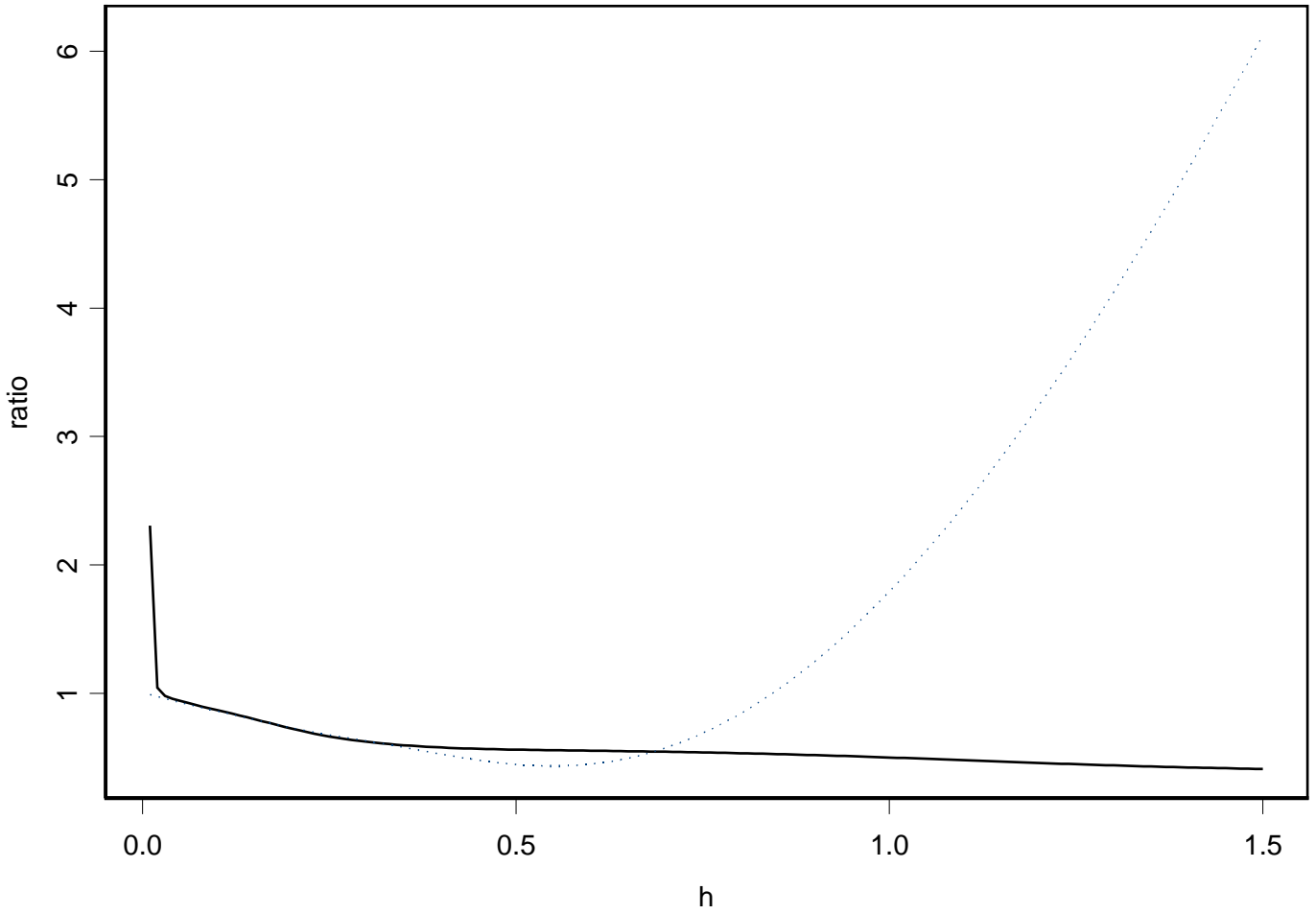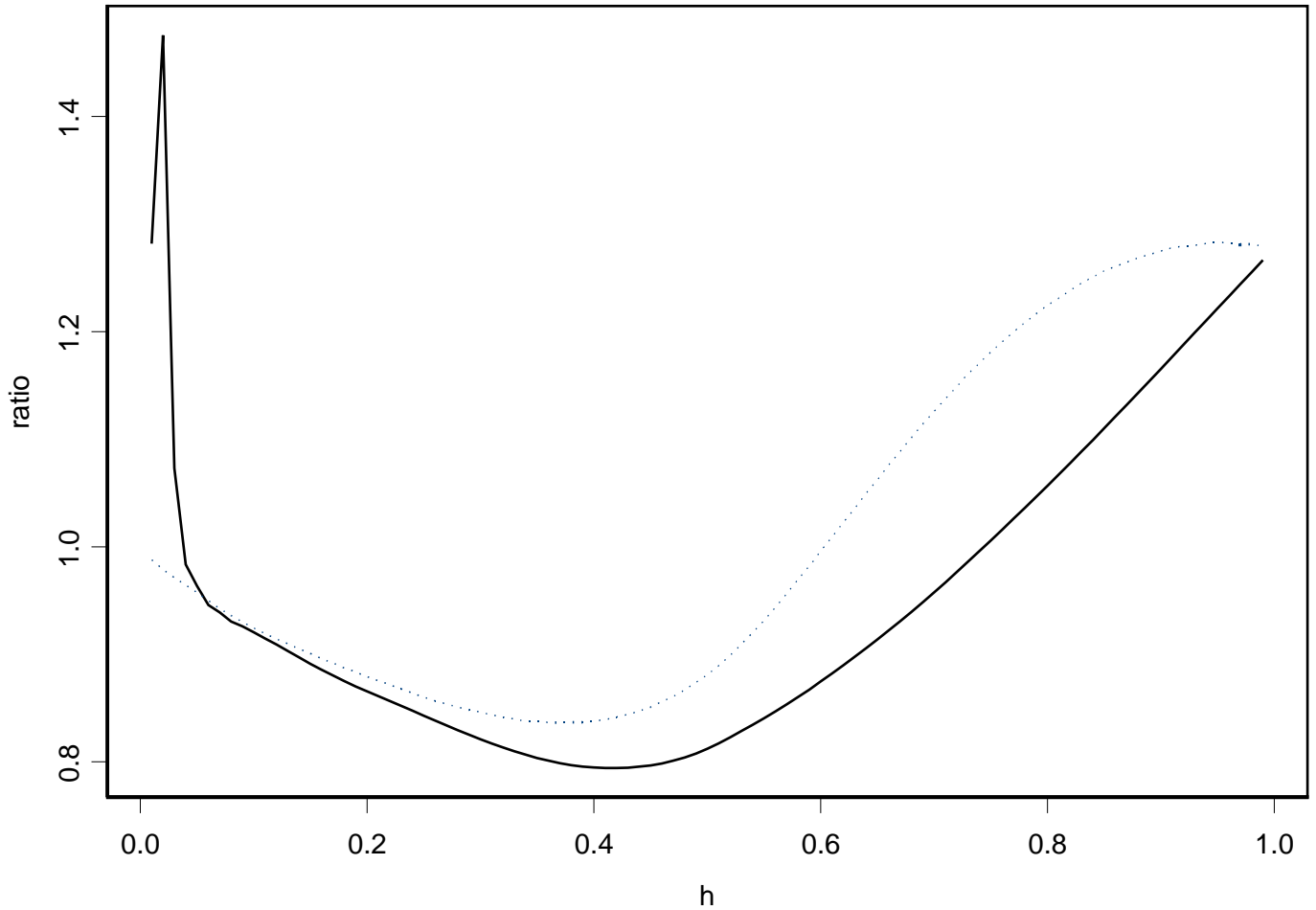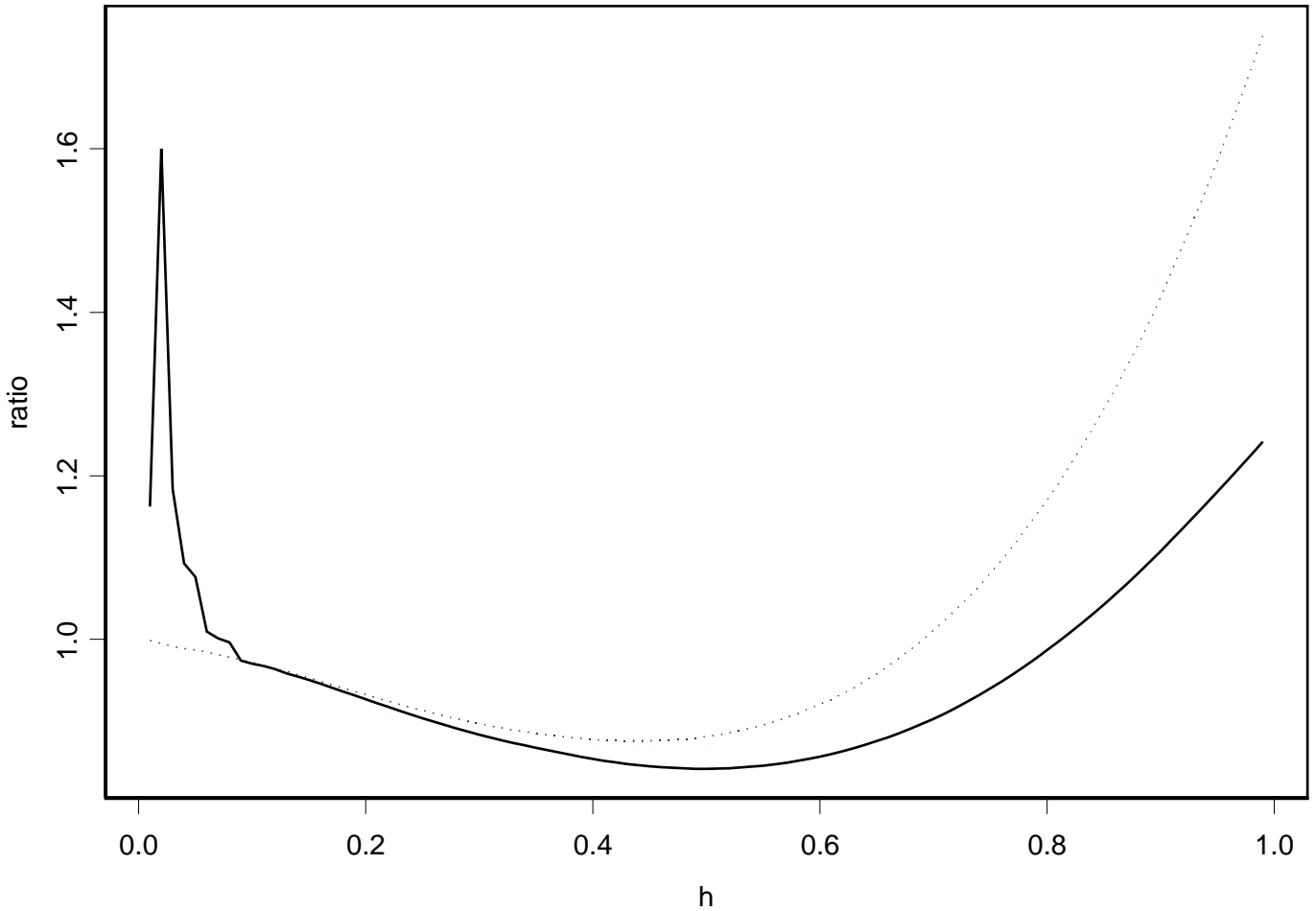
Figure 1: *The solid line and dotted line represent the ratio of the mean squared error of local linear estimator $\hat{F}_n(x)$ to that of product limit estimator $F_n(x)$ and the ratio of the mean squared error of kernel smooth estimator $\bar{F}_n(x)$ to that of product limit estimator $F_n(x)$, respectively. We took $\alpha = 3/7$.*
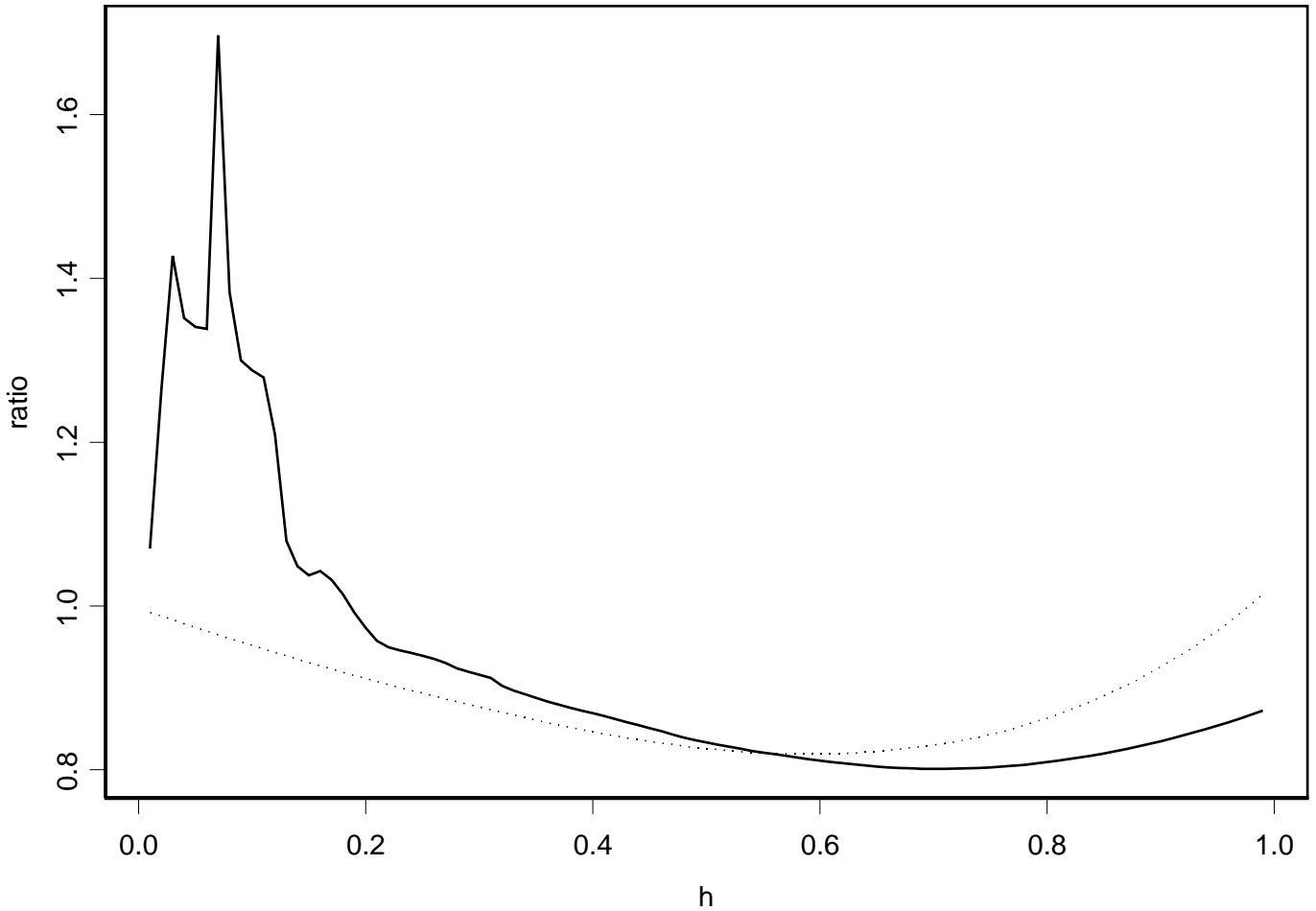
Figure 2: *The solid line and dotted line represent the ratio of the mean squared error of local linear estimator $\hat{F}_n(x)$ to that of product limit estimator $F_n(x)$ and the ratio of the mean squared error of kernel smooth estimator $\bar{F}_n(x)$ to that of product limit estimator $F_n(x)$, respectively. We took $\alpha = 3/7$.*
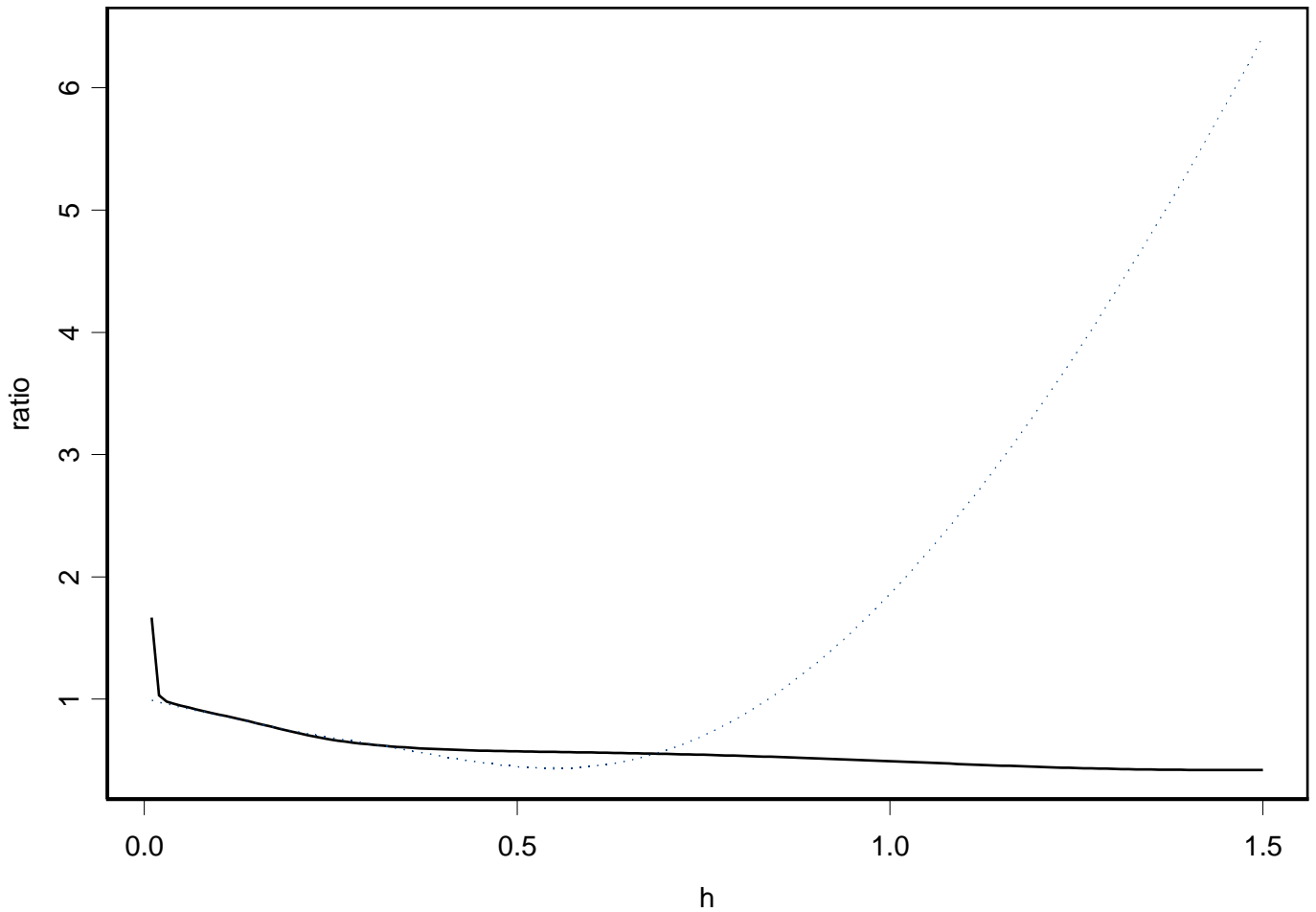
Figure 3: *The solid line and dotted line represent the ratio of the mean squared error of local linear estimator $\hat{F}_n(x)$ to that of product limit estimator $F_n(x)$ and the ratio of the mean squared error of kernel smooth estimator $\bar{F}_n(x)$ to that of product limit estimator $F_n(x)$, respectively. We took $\alpha = 3/7$.*

Figure 4: *The solid line and dotted line represent the ratio of the mean squared error of local linear estimator $\hat{F}_n(x)$ to that of product limit estimator $F_n(x)$ and the ratio of the mean squared error of kernel smooth estimator $\bar{F}_n(x)$ to that of product limit estimator $F_n(x)$, respectively. We took $\alpha = 3/7$.*
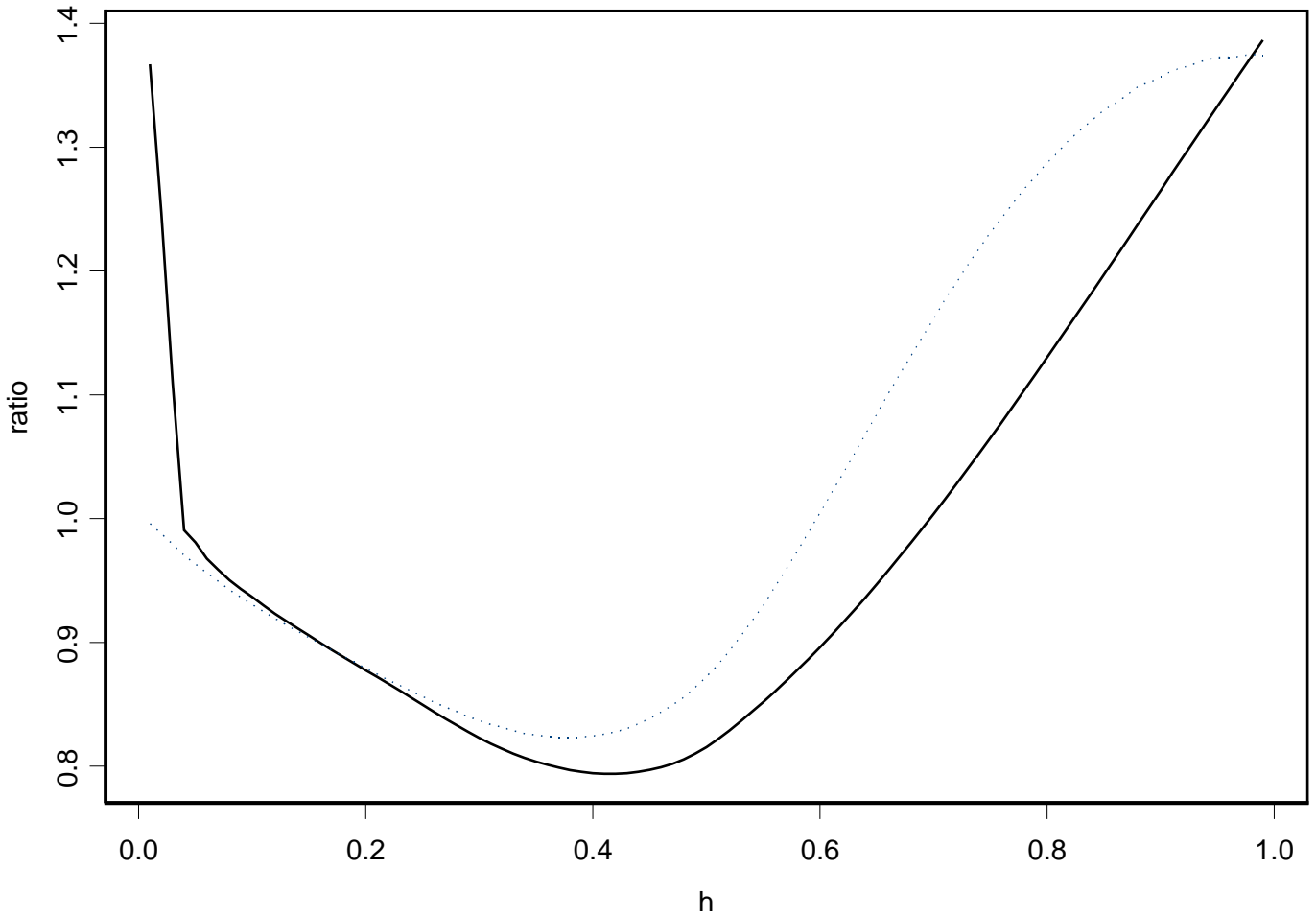
## x=-log(0.8)



Figure 5: *The solid line and dotted line represent the ratio of the mean squared error of local linear estimator $\hat{F}_n(x)$ to that of product limit estimator $F_n(x)$ and the ratio of the mean squared error of kernel smooth estimator $\bar{F}_n(x)$ to that of product limit estimator $F_n(x)$, respectively. We took $\alpha = 1/19$.*
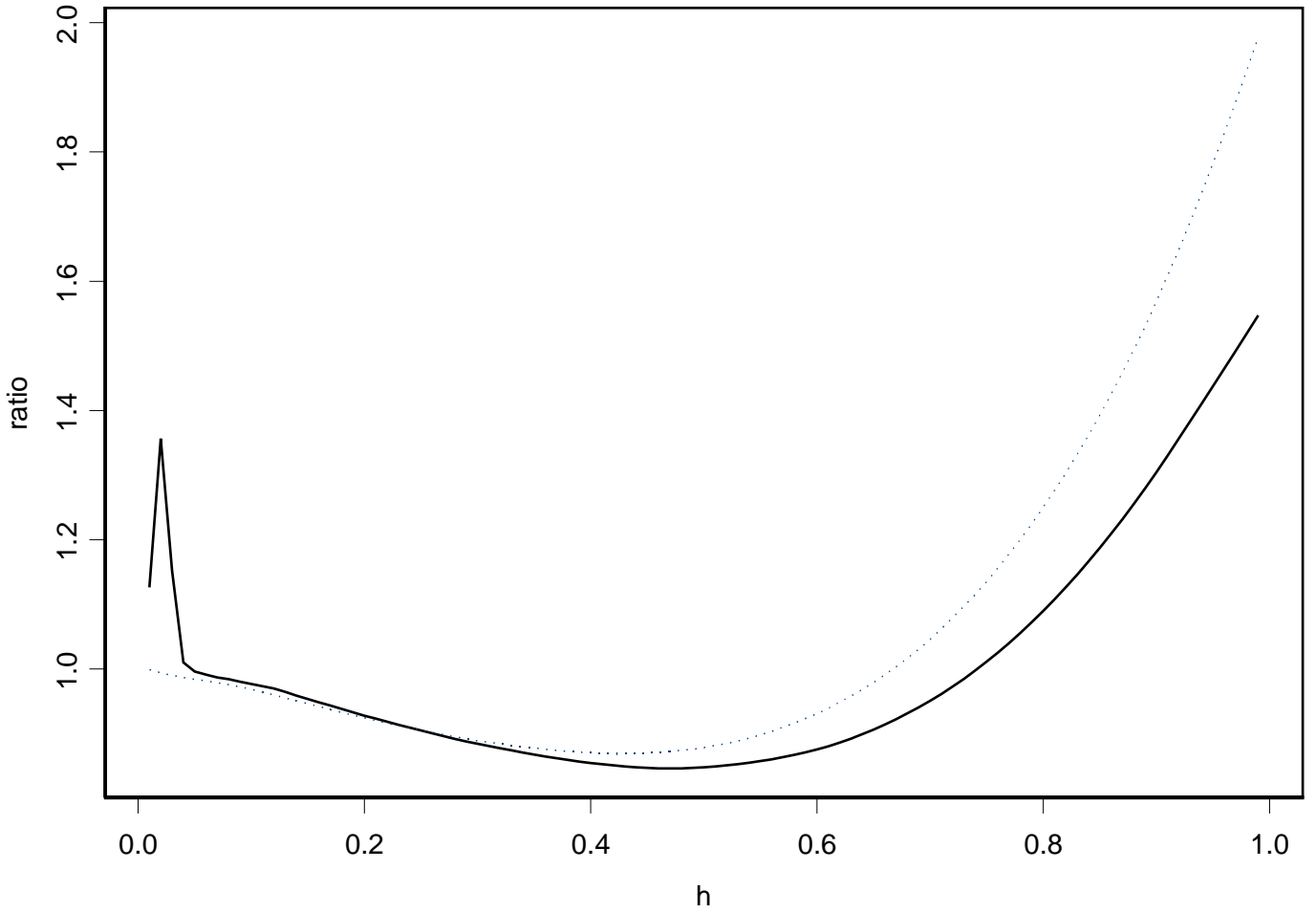
Figure 6: *The solid line and dotted line represent the ratio of the mean squared error of local linear estimator $\hat{F}_n(x)$ to that of product limit estimator $F_n(x)$ and the ratio of the mean squared error of kernel smooth estimator $\bar{F}_n(x)$ to that of product limit estimator $F_n(x)$, respectively. We took $\alpha = 1/19$.*

Figure 7: *The solid line and dotted line represent the ratio of the mean squared error of local linear estimator $\hat{F}_n(x)$ to that of product limit estimator $F_n(x)$ and the ratio of the mean squared error of kernel smooth estimator $\bar{F}_n(x)$ to that of product limit estimator $F_n(x)$, respectively. We took $\alpha = 1/19$.*
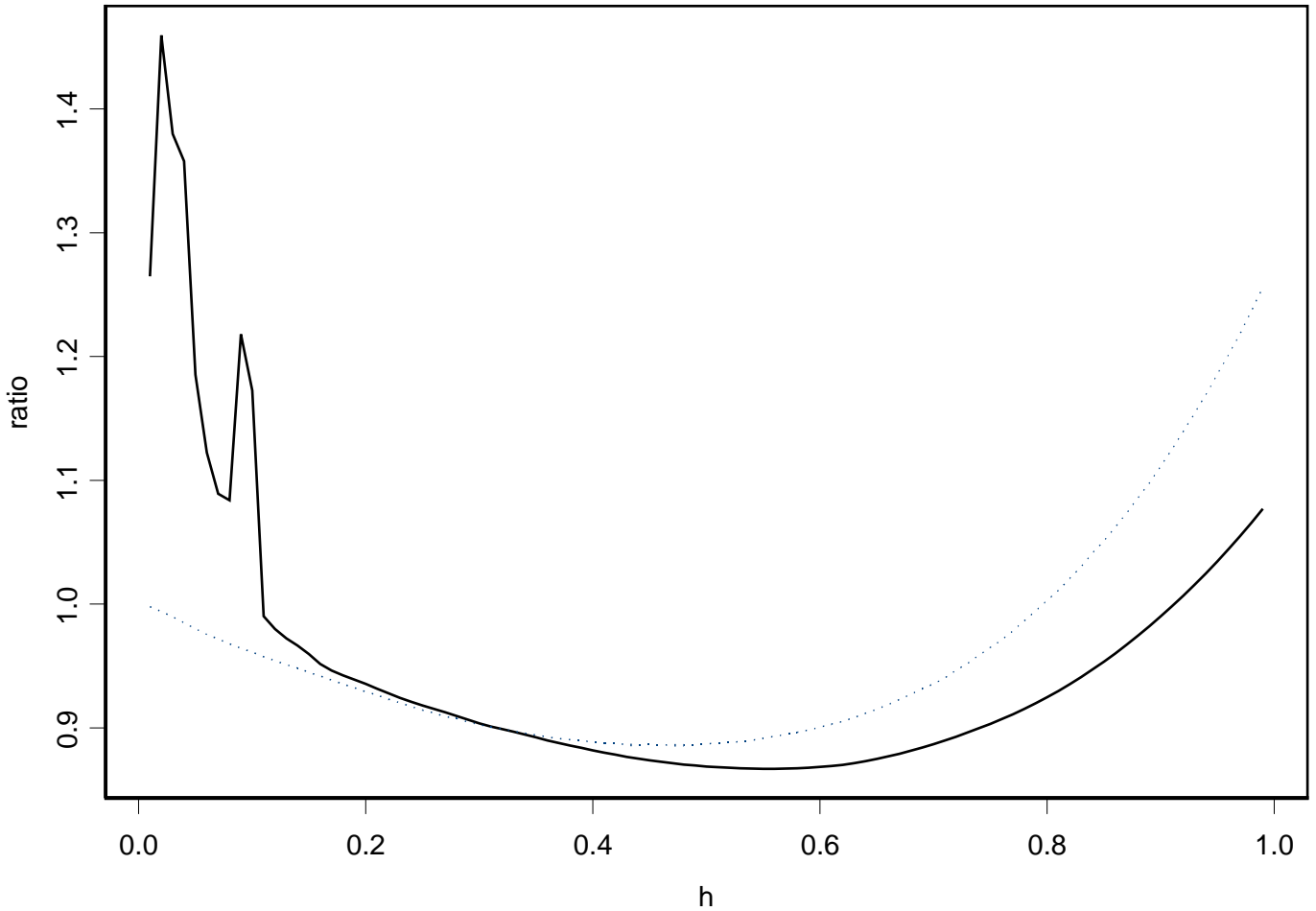
Figure 8: *The solid line and dotted line represent the ratio of the mean squared error of local linear estimator $\hat{F}_n(x)$ to that of product limit estimator $F_n(x)$ and the ratio of the mean squared error of kernel smooth estimator $\bar{F}_n(x)$ to that of product limit estimator $F_n(x)$, respectively. We took $\alpha = 1/19$.*